

Parallelization of the Tau-Leap Coarse-Grained Monte Carlo Method for Efficient and Accurate Simulations on GPUs

Lifan Xu¹, Stuart Collins², Dionisios G. Vlachos², and Michela Taufer¹ ¹Global Computing Lab, Computer and Info. Sciences, University of Delaware ²Chemical Engineering, University of Delaware



Abstract

Recent efforts have outlined the potential of GPUs for Monte Carlo scientific applications. In this poster, we contribute to this effort by exploring the GPU potential for the tau-leaping Coarse-Grained Monte Carlo (CGMC) method. CGMC simulations are important tools for studying phenomena such as catalysis and crystal growth. Existing parallelization of other CGMC method do not support very large molecular system simulations. Our new algorithm on GPUs provides scientists with a much faster way to study very large molecular systems (faster than on traditional HPC clusters) with the same accuracy.

The efficient parallelization of the tau-leaping method for GPUs includes the redesign of both the algorithm and its data structures to address the significantly different GPU memory organization and the GPU multi-thread programming paradigm. The poster describes the parallelization process and the algorithm performance, scalability, and accuracy. To our knowledge, this is the first implementation of this algorithm for GPUs.

• Group neighboring microscopic sites together into "coarse-grained" cells

Coarse Grained Kinetic Monte Carlo Model (CGMC)

• Apply a closure at the stochastic level to resident molecules to describe their distribution in the cell



In the simplest closure above, molecules within cells are assumed to be well mixed. The molecules of each cell are allowed to interact with and diffuse to nearby cells.

CGMC Algorithm on GPUs



Implementation using global memory One thread simulates events in one cell

Multi-thread Implementation of CGMC

Performance



• Use 2-layer algorithm

Thread



Multi-GPU Implementation

Pseudo-code of multi-GPU implementation:

- CPU gets number of GPU and assigns simulation cells to GPUs
- CPU packs data and sends it to each GPU
- Each GPU copies data to its own memory and calls kernel function for two leaps
- All GPUs copy new data back to CPU
- CPU sends new data to GPUs and start next leap

Configuration:

- Multiple GPUs + OpenMP with portable pinned memory (ppm), mapped pinned memory (mpm), and write combined memory (wcm)
- Multiple GPUs + OpenMP with portable pinned memory (ppm), mapped pinned memory (mpm)
- Multiple GPUs + OpenMP with portable pinned memory (ppm)

Conclusions and Future Work

Our contribution:

- Provide scientists with a much faster tau-leaping algorithm to study molecular systems with the same accuracy
- Prove that GPUs have tremendous computational horsepower
- Identify the most suitable level of parallelism for different molecular system sizes

Performance:

• Our tau-leaping algorithm is more than 100 times faster than the sequential version on CPU

Accuracy:

• GPU results are exactly the same as on CPU

Future Work:

coalescing

• Improve performance of multi-GPUs with portable pinned memory

Acknowledgments



BLOCK



A, B, and C

16

32

64

128 256

384

512 1024 1536 2048 4096 8192 12288





• Number of events per msec is 120X faster than on a single CPU

• Max number of cells grows from 32,768 on single GPU to 102,400 on



